

Gradient descent for neural network.

Parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$
 $(n^{[0]}, n^{[0]}) (n^{[1]}, 1) (n^{[2]}, n^{[1]}) (n^{[2]}, 1)$

Cost function.: $J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i)$.

→ Gradient Descent:

Repeat:

Compute predicts: $(\hat{y}^{(0)}, \dots, \hat{y}^{(m)})$.

$d w^{[1]} = \frac{\partial J}{\partial w^{[1]}}$, $d b^{[1]} = \frac{\partial J}{\partial b^{[1]}}$, ... other partial derivatives.

$$w^{[1]} := w^{[1]} - \alpha d w^{[1]}$$

$$b^{[1]} := b^{[1]} - \alpha d b^{[1]}$$

$$w^{[2]} := w^{[2]} - \alpha d w^{[2]}$$

... update all parameters using their partial derivatives and the learning rate.

Mathematical Formulas:

forward propagation:

$$z^{[1]} = w^{[1]} x + b^{[1]}$$

$$A^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = w^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(z^{[2]})$$

These operations are all vectorized.

⋮

Back Propagation :-

Rewind about Computing gradients for logistic Regression:-

$$\begin{array}{l} x \\ w \\ b \end{array} \rightarrow z = w^T x + b \rightarrow a = \sigma(z) \rightarrow L(a, y)$$
$$dw = dz \cdot x$$
$$db = dz \cdot a' = da \cdot g'(z) = a - y$$
$$dz = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} = -\frac{1}{a} + \frac{1-y}{1-a}$$

Back Propagation for Neural Network :-

$$\begin{array}{l} x \\ w^{[1]} \\ b^{[1]} \end{array} \rightarrow z^{[1]} = w^{[1]} x + b^{[1]} \rightarrow a^{[1]} = \sigma(z^{[1]})$$
$$dz^{[1]} = W^{[2]} dZ^{[2]}$$
$$dW^{[1]} = dz^{[1]} x^T$$
$$db^{[1]} = dz^{[1]} a^{[1]T}$$
$$w^{[2]} \\ b^{[2]} \end{array} \rightarrow z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$
$$dZ^{[2]} = a^{[2]} y$$
$$dW^{[2]} = dZ^{[2]} a^{[1]T}$$
$$db^{[2]} = dZ^{[2]}$$
$$da^{[2]} = \frac{\partial L}{\partial a}$$
$$L(a, y)$$

Summary of gradient descent.

$$dz^{[2]} = a^{[2]} - y$$

$$dW^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = W^{[2]T} dz^{[2]} + g'(z^{[1]})$$

$$dW^{[1]} = dz^{[1]} x^T$$

$$db^{[1]} = dz^{[1]}$$

Loss function is only used for prediction (last layer).

Vectorize Backward propagation:

$$dz^{[2]} = A^{[2]} - Y$$

$$dw^{[2]} = \frac{1}{m} dz^{[2]} A^{[1] T}$$

$$db^{[2]} = \frac{1}{m} np.sum(dz^{[2]}, axis=1, keepdims=True)$$

$$dz^{[1]} = w^{[2] T} dz^{[2]} * g^{[1]'}(\bar{z}^{[1]})$$

$$dw^{[1]} = \frac{1}{m} dz^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} np.sum(dz^{[1]}, axis=1, keepdims=True).$$